

General Disclaimer

One or more of the Following Statements may affect this Document

- This document has been reproduced from the best copy furnished by the organizational source. It is being released in the interest of making available as much information as possible.
- This document may contain data, which exceeds the sheet parameters. It was furnished in this condition by the organizational source and is the best copy available.
- This document may contain tone-on-tone or color graphs, charts and/or pictures, which have been reproduced in black and white.
- This document is paginated as submitted by the original source.
- Portions of this document are not fully legible due to the historical nature of some of the material. However, it is the best reproduction available from the original submission.



DEPARTMENT OF MATHEMATICS

UNIVERSITY OF HOUSTON

HOUSTON, TEXAS

NASA CR-

147521

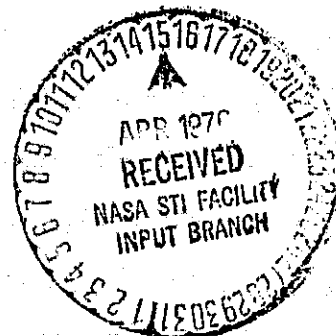
(NASA-CR-147521) AN ITERATIVE PROCEDURE FOR
OBTAINING MAXIMUM-LIKELIHOOD ESTIMATES OF
THE PARAMETERS FOR A MIXTURE OF NORMAL
DISTRIBUTIONS, ADDENDUM (Houston Univ.)
11 p HC \$3.50

N76-20914

Unclas
21495

CSSL 12A G3/65

ADDENDUM TO "AN ITER PROC FOR
OBTAINING MAX-LIKELIHOOD ESTS OF
THE PARAM FOR A MIX OF NOR LISTS
B. C. PETERS, JR & H. F. WALKER
REPORT #47 SEPTEMBER 1975



PREPARED FOR
EARTH OBSERVATION DIVISION, JSC
UNDER
CONTRACT NAS-9-12777

3801 CULLEN BLVD.
HOUSTON, TEXAS 77004

*Addendum to "An Iterative Procedure for
Obtaining Maximum-Likelihood Estimates of
the Parameters for a Mixture of Normal Distributions"*

by

B. Charles Peters, Jr.

*NASA/National Research Council Research Associate
Earth Observations Division, Johnson Space Center*

and

Homer F. Walker

*Department of Mathematics, University of Houston
Houston, Texas*

*Report #47
September, 1975
NAS 9-12777*

Addendum to "An Iterative Procedure for
Obtaining Maximum-Likelihood Estimates of
the Parameters for a Mixture of Normal Distributions"

by

B. Charles Peters, Jr.

NASA/National Research Council Research Associate
Earth Observations Division, Johnson Space Center

and

Homer F. Walker

Department of Mathematics, University of Houston
Houston, Texas

1. Introduction.

In this report, we discuss new results and insights concerning an iterative procedure introduced in [1] for obtaining maximum-likelihood estimates of the parameters for a mixture of normal distributions. For any questions concerning notation, definitions, etc., the reader is referred to that report.

The iterative procedure in question is the following: Beginning with some starting value $\begin{pmatrix} \bar{\alpha}(1) \\ \bar{\mu}(1) \\ \bar{\Sigma}(1) \end{pmatrix}$ in the space $\mathcal{A} \oplus \mathcal{M} \oplus \mathcal{S}$ introduced in [1], define

successive iterates inductively by the relationship

$$(*) \quad \begin{pmatrix} \bar{\alpha}^{(k+1)} \\ \bar{\mu}^{(k+1)} \\ \bar{\Sigma}^{(k+1)} \end{pmatrix} = \mathbb{I}_\epsilon(\bar{\alpha}^{(k)}, \bar{\mu}^{(k)}, \bar{\Sigma}^{(k)})$$

given in [1]. It is shown in [1] that, with probability approaching 1 as the sample size N approaches infinity, this procedure converges locally to the consistent maximum-likelihood estimate whenever ϵ is sufficiently small. (In particular, $\epsilon < \frac{4}{m(n+1)(n+2)}$ guarantees the local convergence of this procedure in probability.)

In this report, we prove that, in probability, the procedure $(*)$ converges locally to the consistent maximum-likelihood estimate whenever $0 < \epsilon < 2$. We also show that the ϵ which yields optimal local convergence rates lies between 1 and 2. In fact, the optimal ϵ is near 1, if the component populations are widely separated, and near 2 if the component populations have nearly identical means and covariance matrices.

1. Local Convergence.

As in [1], we say that \mathbb{I}_ϵ is locally contractive (in a norm $\|\cdot\|$ on $\mathcal{A} \oplus \mathcal{M} \oplus \mathcal{S}$) near $\begin{pmatrix} \bar{\alpha} \\ \bar{\mu} \\ \bar{\Sigma} \end{pmatrix} \in \mathcal{A} \oplus \mathcal{M} \oplus \mathcal{S}$ if there is a number λ , $0 \leq \lambda < 1$

such that

$$\left\| \mathbb{I}_\epsilon(\bar{\alpha}, \bar{\mu}, \bar{\Sigma}) - \begin{pmatrix} \bar{\alpha} \\ \bar{\mu} \\ \bar{\Sigma} \end{pmatrix} \right\| \leq \lambda \left\| \begin{pmatrix} \bar{\alpha}' \\ \bar{\mu}' \\ \bar{\Sigma}' \end{pmatrix} - \begin{pmatrix} \bar{\alpha} \\ \bar{\mu} \\ \bar{\Sigma} \end{pmatrix} \right\|$$

whenever $\begin{pmatrix} \bar{\alpha}' \\ \bar{\mu}' \\ \bar{\Sigma}' \end{pmatrix}$ lies sufficiently near $\begin{pmatrix} \bar{\alpha} \\ \bar{\mu} \\ \bar{\Sigma} \end{pmatrix}$. Our result is the following.

Theorem. With probability approaching 1 as N approaches infinity, \mathbb{I}_ϵ is a locally contractive operator (in a norm to be defined on $\mathcal{O} \oplus \mathcal{M} \oplus \mathcal{S}$) near the consistent maximum-likelihood estimate whenever $0 < \epsilon < 2$.

Corollary. With probability approaching 1 as N approaches infinity, the iterative procedure (*) converges locally to the consistent maximum-likelihood estimate whenever $0 < \epsilon < 2$.

Proof: As observed in [1], the theorem will be proved if it can be shown that, for $0 < \epsilon < 2$, $E(\nabla \mathbb{I}_\epsilon(\bar{\alpha}^0, \bar{\mu}^0, \bar{\Sigma}^0))$ has operator norm less than 1 with respect to some vector norm on $\mathcal{O} \oplus \mathcal{M} \oplus \mathcal{S}$. (Throughout this note, the superscript "0" indicates that the superscripted parameters are the true parameters of the mixture density.) For $i=1, \dots, m$, let $\langle \cdot, \cdot \rangle_1'$ and $\langle \cdot, \cdot \rangle_1''$ be the inner products on R^n and the space of real, symmetric $n \times n$ matrices introduced in [1], i.e., let

$$\langle v, w \rangle_1' = v^T (\alpha_1^0 \Sigma_1^0)^{-1} w \quad \text{for } v, w \in R^n,$$

$$\langle A, B \rangle_1'' = \text{tr} \left\{ A \left(-\frac{\alpha_1^0}{2} \Sigma_1^0 \right)^{-1} B^T \right\} \quad \text{for real, symmetric } n \times n \text{ } A \text{ and } B.$$

These inner products, together with scalar multiplication on R^1 , induce an inner product $\langle \cdot, \cdot \rangle$ on $\mathcal{A} \oplus \mathcal{M} \oplus \mathcal{S}$. Now $E(\nabla \tilde{\Phi}_\epsilon(\bar{\alpha}^\circ, \bar{\mu}^\circ, \bar{\Sigma}^\circ)) = I = \epsilon QR$, where

$$Q = \begin{pmatrix} (\text{diag } \alpha_i^\circ) & 0 & 0 \\ 0 & I & 0 \\ 0 & 0 & (\text{diag } \Sigma_i^\circ) \end{pmatrix}$$

and

$$\int_{R_n} \begin{pmatrix} \frac{p_1}{p} \\ \vdots \\ \frac{p_m}{p} \\ \frac{p_1}{p}(x-\mu_1^\circ) \\ \vdots \\ \frac{p_m}{p}(x-\mu_m^\circ) \\ \frac{p_1}{p}[\Sigma_1^{\circ-1}(x-\mu_1^\circ)(x-\mu_1^\circ)^T - I] \\ \vdots \\ \frac{p_m}{p}[\Sigma_m^{\circ-1}(x-\mu_m^\circ)(x-\mu_m^\circ)^T - I] \end{pmatrix} \langle \begin{pmatrix} \frac{p_1}{p} \\ \vdots \\ \frac{p_m}{p} \\ \frac{p_1}{p}(x-\mu_1^\circ) \\ \vdots \\ \frac{p_m}{p}(x-\mu_m^\circ) \\ \frac{p_1}{p}[\Sigma_1^{\circ-1}(x-\mu_1^\circ)(x-\mu_1^\circ)^T - I] \\ \vdots \\ \frac{p_m}{p}[\Sigma_m^{\circ-1}(x-\mu_m^\circ)(x-\mu_m^\circ)^T - I] \end{pmatrix} \rangle_p dx$$

One sees that the theorem will be proved if it can be shown that, with respect to some vector norm on $\mathcal{A} \oplus \mathcal{M} \oplus \mathcal{S}$, the operator norm of QR is no greater than 1. Since QR is positive definite and symmetric with respect to the inner product $\langle \cdot, Q^{-1} \cdot \rangle$, it follows that the theorem will be proved if it can be shown that $\langle V, Q^{-1}[QR]V \rangle = \langle V, RV \rangle \leq \langle V, Q^{-1}V \rangle$ for $V \in \mathcal{A} \oplus \mathcal{M} \oplus \mathcal{S}$.

For

$$V = \begin{pmatrix} y_1 \\ \vdots \\ y_m \\ v_1 \\ \vdots \\ v_m \\ A_1 \\ \vdots \\ A_m \end{pmatrix} \in \mathcal{A} \otimes \mathcal{M} \otimes \mathcal{S},$$

one has

$$\begin{aligned} \langle V, RV \rangle &= \int_{\mathbb{R}^n} \left(\sum_{i=1}^m y_i \frac{p_i}{p} + \sum_{i=1}^m v_i^T (\alpha_i^\circ \Sigma_i^{\circ-1}) \frac{p_i}{p} (x - \mu_i^\circ) + \right. \\ &\quad \left. \sum_{i=1}^m \text{tr} \left(A_i \left(\frac{1}{2} \Sigma_i^{\circ-1} \right) \frac{p_i}{p} [\Sigma_i^{\circ-1} (x - \mu_i^\circ) (x - \mu_i^\circ)^T - I] \right) \right)^2 p \, dx \\ &= \int_{\mathbb{R}^n} \left(\sum_{i=1}^m [\alpha_i^{\circ-1} y_i + v_i^T \Sigma_i^{\circ-1} (x - \mu_i^\circ) + \text{tr} \{ A_i (\frac{1}{2} \Sigma_i^{\circ-1}) [\Sigma_i^{\circ-1} (x - \mu_i^\circ) (x - \mu_i^\circ)^T - I] \}] \frac{\alpha_i^\circ p_i}{p} \right)^2 p \, dx \\ &\leq \int_{\mathbb{R}^n} \left(\sum_{i=1}^m [\alpha_i^{\circ-1} y_i + v_i^T \Sigma_i^{\circ-1} (x - \mu_i^\circ) + \text{tr} \{ A_i (\frac{1}{2} \Sigma_i^{\circ-1}) [\Sigma_i^{\circ-1} (x - \mu_i^\circ) (x - \mu_i^\circ)^T - I] \}]^2 \frac{\alpha_i^\circ p_i}{p} \right) p \, dx \end{aligned}$$

by Schwarz's inequality. If the squared expressions in the last sum above are written out in full, one sees that the integrals of the cross terms in these expressions vanish. Consequently,

$$\langle V, RV \rangle \leq \int_{\mathbb{R}^n} \left(\sum_{i=1}^m [\alpha_i^{\circ-2} y_i^2 + (v_i^T \Sigma_i^{\circ-1} (x - \mu_i^\circ))^2 + (\text{tr} \{ A_i (\frac{1}{2} \Sigma_i^{\circ-1}) [\Sigma_i^{\circ-1} (x - \mu_i^\circ) (x - \mu_i^\circ)^T - I] \})^2 \alpha_i^\circ p_i] \right) p \, dx$$

Now

$$(1) \quad \int_{R^n} \alpha_1^{\circ-1} y_1^2 p_1 dx = \alpha_1^{\circ-1} y_1^2$$

$$(2) \quad \int_{R^n} (v_1^T \Sigma_1^{\circ-1} (x - \mu_1^{\circ}))^2 \alpha_1^{\circ} p_1 dx = \int_{R^n} v_1^T \Sigma_1^{\circ-1} (x - \mu_1^{\circ}) (x - \mu_1^{\circ})^T \Sigma_1^{\circ-1} v_1 \alpha_1^{\circ} p_1 dx$$

$$= \langle v_1, v_1 \rangle_1$$

$$(3) \quad \int_{R^n} (\text{tr}\{A_1 (\frac{1}{2} \Sigma_1^{\circ-1}) [\Sigma_1^{\circ-1} (x - \mu_1^{\circ}) (x - \mu_1^{\circ})^T - I]\}^2 \alpha_1^{\circ} p_1 dx = \langle A_1, \Sigma_1^{\circ-1} A_1 \rangle_1''$$

(A proof of (3) follows below.) From (1), (2), and (3), one concludes that

$$\langle V, RV \rangle \leq \sum_{i=1}^m \alpha_1^{\circ-1} y_i^2 + \langle v_1, v_1 \rangle_1 + \langle A_1, \Sigma_1^{\circ-1} A_1 \rangle_1'' = \langle V, Q^{-1} V \rangle.$$

This completes the proof of the theorem.

Proof of (3): Setting $y = \Sigma_1^{\circ-1/2} (x - \mu_1^{\circ})$ and

$$I = \int_{R^n} (\text{tr}\{A_1 (\frac{1}{2} \Sigma_1^{\circ-1}) [\Sigma_1^{\circ-1} (x - \mu_1^{\circ}) (x - \mu_1^{\circ})^T - I]\}^2 \alpha_1^{\circ} p_1 dx,$$

one obtains

$$I = \frac{\alpha_1^{\circ}}{4} \int_{R^n} (\text{tr}\{A_1 [\Sigma_1^{\circ-1/2} y y^T \Sigma_1^{\circ-1/2} - \Sigma_1^{\circ-1}]\}^2 p_0 dy,$$

where $p_0 \sim N(0, I)$. Denoting $\Sigma_1^{\circ-1/2} A_1 \Sigma_1^{\circ-1/2} = B = (b_{jk})$,

one then derives

$$\begin{aligned}
 I &= \frac{\alpha_1^0}{4} \int_{R^n} (\text{tr}\{B[yy^T - I]\})^2 p_0 dy \\
 &= \frac{\alpha_1^0}{4} \int_{R^n} [(\text{tr}\{Byy^T\})^2 - 2\text{tr}\{B\}\text{tr}\{Byy^T\} + (\text{tr}\{B\})^2] p_0 dy \\
 &= \frac{\alpha_1^0}{4} \left\{ \sum_{j,k,p,q} \beta_{jk} \beta_{pq} \int_{R^n} y_k y_j y_q y_p p_0 dy - 2 \text{tr}\{B\} \sum_{j,k} \beta_{jk} \int_{R^n} y_k y_j p_0 dy + (\text{tr}\{B\})^2 \right\} \\
 &= \frac{\alpha_1^0}{4} \left\{ \sum_k \sum_{p \neq k} \beta_{kk} \beta_{pp} + \sum_k \sum_{j \neq k} \beta_{jk} \beta_{jk} + \sum_k \sum_{j \neq k} \beta_{jk} \beta_{kj} + 3 \sum_k \beta_{kk}^2 - 2(\text{tr}\{B\})^2 + (\text{tr}\{B\})^2 \right\} \\
 &= \frac{\alpha_1^0}{2} \text{tr}\{B^2\} = \frac{\alpha_1^0}{2} \text{tr}\{\Sigma_1^{-1/2} A_1 \Sigma_1^{-1} A_1 \Sigma_1^{-1/2}\} = \text{tr}\{A_1 (\frac{\alpha_1^0}{2} \Sigma_1^{-1}) (\Sigma_1^{-1} A_1)^T\} \\
 &= \langle A_1, \Sigma_1^{-1} A_1 \rangle_1''
 \end{aligned}$$

3. The optimal ϵ .

From the proof of the theorem, one sees that, asymptotically as N approaches infinity, the value of ϵ which yields optimal local convergence rates is that which minimizes the spectral radius of $E(\nabla \bar{\Phi}_\epsilon(\bar{\alpha}^o, \bar{\mu}^o, \bar{\Sigma}^o))$. (Indeed, $E(\nabla \bar{\Phi}_\epsilon(\bar{\alpha}^o, \bar{\mu}^o, \bar{\Sigma}^o)) = I - \epsilon QR$ is symmetric with respect to the inner product $\langle \cdot, Q^{-1} \cdot \rangle$; hence, its operator norm with respect to this inner product is equal to its spectral radius.) Letting ρ and τ denote, respectively, the largest and smallest eigenvalues of QR , one verifies that the spectral radius of $E(\nabla \bar{\Phi}_\epsilon(\bar{\alpha}^o, \bar{\mu}^o, \bar{\Sigma}^o))$ is minimized when $1 - \epsilon \tau = \epsilon \rho - 1$, i.e., when $\epsilon = \frac{2}{\rho + \tau}$. Now $\rho = 1$ always, for it follows from the proof of the theorem that ρ is

never greater than 1, and

$$\begin{pmatrix} \alpha_1^0 \\ \vdots \\ \alpha_m^0 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \in \mathcal{O} \oplus \mathcal{N} \oplus \mathcal{S}$$

is always an eigenvector of QR with eigenvalue 1. Thus optimal convergence rates are obtained when $\epsilon = \frac{2}{1+\tau}$, where τ lies between 0 and 1. In particular, the best choice of ϵ lies between 1 and 2.

Suppose that the component populations in the mixture are "widely separated" in the sense that each pair (μ_i^0, Σ_i^0) differs greatly from every other such pair. Then

$$\left(\frac{\alpha_i^0 p_i(x)}{p(x)} \right) \left(\frac{\alpha_j^0 p_j(x)}{p(x)} \right) \approx \delta_{ij} \quad \text{for } x \in \mathbb{R}^n \text{ and } i, j = 1, \dots, m,$$

and one verifies that $QR \approx I$. Consequently, optimal convergence rates are obtained for an ϵ near 1 and, for the optimal ϵ , $E(\nabla \bar{\Phi}_\epsilon(\bar{\alpha}^0, \bar{\mu}^0, \bar{\Sigma}^0)) = I - \epsilon QR \approx 0$. Thus for mixtures whose component populations are "widely separated", optimal convergence rates are obtained for an ϵ near 1, and rapid first-order convergence can be expected for this ϵ .

Now suppose that the component populations in the mixture are such that each pair (μ_i^0, Σ_i^0) differs little from every other such pair. Then

$p(x) \approx p_i(x)$ and $\frac{p_i(x)}{p(x)} \approx 1$ for $x \in R^n$ and $i = 1, \dots, m$, and one verifies that the smallest eigenvalue of QR is near zero. It follows that optimal convergence rates are obtained for an ϵ near 2. In this case, the spectral radius of $E(\nabla \bar{J}_\epsilon(\bar{\alpha}^0, \bar{\mu}^0, \bar{\Sigma}^0))$ is near 1, even for the optimal value of ϵ ; hence, slow first-order convergence is to be expected.

We conclude by observing that the major practical implication of this note is that the iterative procedure under consideration converges whenever the step-size ϵ lies in an interval which is completely independent of the particular mixture problem at hand. It is readily ascertained that this cannot be said for the regular steepest descent procedure

$$\begin{aligned}\alpha_1^{(q+1)} &= \alpha_1^{(q)} + \epsilon \left[\frac{1}{N} \sum_{k=1}^N \frac{p_1(x_k)}{p(x_k)} - \frac{1}{mN} \sum_{j=1}^m \sum_{k=1}^N \frac{p_j(x_k)}{p(x_k)} \right] \\ \mu_1^{(q+1)} &= \mu_1^{(q)} + \epsilon \left[\frac{1}{N} \sum_{k=1}^N \frac{\alpha_1^{(q)} p_1(x_k)}{p(x_k)} \Sigma_1^{(q)-1} (x_k - \mu_1^{(q)}) \right] \\ \Sigma_1^{(q+1)} &= \Sigma_1^{(q)} + \epsilon \left[\frac{1}{2N} \sum_{k=1}^N \frac{\alpha_1^{(q)} p_1(x_k)}{p(x_k)} [-\Sigma_1^{(q)-1} + \Sigma_1^{(q)-1} (x_k - \mu_1^{(q)}) (x_k - \mu_1^{(q)})^T \Sigma_1^{(q)-1}] \right].\end{aligned}$$

Thus the procedure considered here offers considerable practical advantages over the steepest descent procedure, even though it is itself a generalized steepest descent (deflected gradient) procedure.

REFERENCE

1. B.C. Peters and H.F. Walker, "An iterative procedure for obtaining maximum-likelihood estimates of the parameters for a mixture of normal distributions, "Report #43, NASA Contract NAS-9-12777, University of Houston, Department of Mathematics.